

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-13959

(43) 公開日 平成7年(1995)1月17日

(51) Int.Cl. <sup>6</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 17/00		7055-2 J		
G 0 1 N 33/68				
G 0 6 F 17/30		8724-5 L	G 0 6 F 15/ 20	D
		9194-5 L	15/ 40	5 3 0 S
審査請求 有 請求項の数 6 O L (全 9 頁)				

(21) 出願番号 特願平4-124817

(22) 出願日 平成4年(1992)5月18日

(71) 出願人 000004237

日本電気株式会社  
東京都港区芝五丁目7番1号

(72) 発明者 馬見塚 拓

東京都港区芝五丁目7番1号日本電気株式  
会社内

(72) 発明者 山西 健司

東京都港区芝五丁目7番1号日本電気株式  
会社内

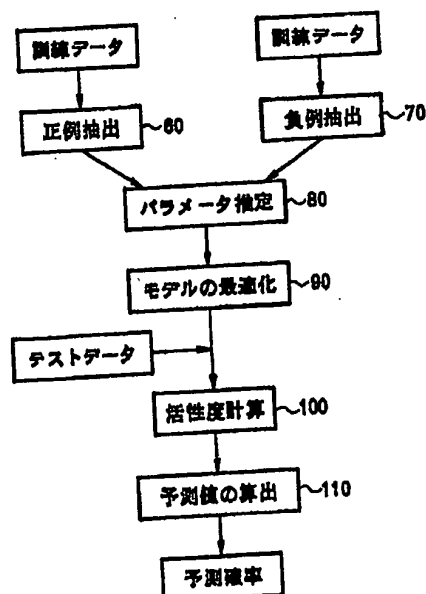
(74) 代理人 弁理士 京本 直樹 (外2名)

(54) 【発明の名称】 タンパク質立体構造予測方法

(57) 【要約】

【目的】 構造未知のタンパク質のアミノ酸配列情報から、それに対応する二次構造 ( $\alpha$ ヘリックス) を高い信頼性で予測する。

【構成】 ステップ60で構造既知及び未知のタンパク質アミノ酸配列を入力とし、それらのアライメント (整合) から二次構造領域の正例を出力し、ステップ70で構造既知のタンパク質アミノ酸配列を入力とし、それらのアライメントから二次構造領域の負例を出力し、ステップ80で前記正例と負例を入力とし、これら訓練データのアミノ酸の実数値属性から確率的規則の実数値パラメータの推定値を出力し、ステップ90で情報量規準を用いて前記確率的規則の最適化したセル数を出力し、ステップ100でテストデータ配列を入力とし、前記確率的規則を使用し、計算を行ったテストデータ配列の各領域に対する活性度を出力し、ステップ110で前記活性度を入力とし、その中から最適値を出力する。



## 【特許請求の範囲】

【請求項1】 タンパク質のアミノ酸配列からタンパク質の構造予測を行うための訓練データを抽出する訓練データ抽出手段と、訓練データから確率的規則の学習を行う学習手段と、学習された確率的規則を用いてテストアミノ酸配列データに対してテストを行うテスト手段とから成ることを特徴とするタンパク質立体構造予測方法。

【請求項2】 前記訓練データ抽出手段が、構造既知のタンパク質のアミノ酸配列に対して、同じファミリーに属するタンパク質のアライメント（整合）をとり、予測対象とする二次構造領域に対応する部分配列を、二次構造領域の正例として抽出するステップと、構造既知のタンパク質の予測対象とする二次構造に対応する部分配列に対して、構造既知のタンパク質からなるデータベースの各配列のアライメントをとり、予測対象とする二次構造に対応しない部分配列を、二次構造領域の負例として抽出するステップとから成ることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項3】 前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の種類から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定することを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項4】 前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の実数値属性から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定することを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項5】 前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の実数値属性から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定するステップと、前記ステップの確率的規則におけるモデルを情報量規準を用いて最適化するステップとから成ることを特徴とする請求項1記載のタンパク質立体構造予測方法。

【請求項6】 前記テストデータ配列に対するテスト手段が、前記学習方法により学習された確率的規則を使用し、テストデータ配列の各領域に対して、その活動度を計算するステップと、計算された活動度の中から最適値を選出するステップとから成ることを特徴とする請求項1記載のタンパク質立体構造予測方法。

## 【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、構造未知のタンパク質アミノ酸配列から、そのタンパク質の立体構造を予測する方法に関する。

【0002】

【従来の技術】 タンパク質のアミノ酸配列情報を用いて、そのタンパク質内の立体構造を予測する問題の一つとして、タンパク質二次構造予測問題がある。二次構造

とは、 $\alpha$ ヘリックスや $\beta$ シートといったタンパク質内部でのまとまりのある構造を指し、二次構造予測問題は、タンパク質のアミノ酸配列情報を用いて、3（あるいは4）種類の二次構造の中から、一次配列の各残基（以下、予測対象となる残基を中心残基とする）に対応する一つの二次構造を予測する問題であり、二次構造予測が可能になることにより、タンパク質の立体的な構造予測も可能になると考えられている。図3は、本発明の二次構造（ $\alpha$ ヘリックス）領域予測方法を示す模式図であるが、従来技術によるタンパク質の二次構造を予測する方法として、例えば、1974年発行の米国の雑誌「バイオケミストリー」（Biochemistry）の第23巻222-245頁掲載のチョウ（Chou）とファスマン（Fasman）による論文「プレディクション オブ プロテイン コンホメーション」（Prediction of protein conformation）（以下、CF法と略す）、1978年発行の米国の雑誌「ジャーナルオブ モレキュラ バイオロジー」（Journal of Molecular Biology）の第120巻97-120頁掲載のガルニエ（Garnier）らによる論文「アナリシス オブ ザ アクキュレシー アンド インプリケーションズ オブ シンプル メソッド フォー プレディクティン グ ザ セコンダリー ストラクチャー オブ グロブular プロテインズ」（Analysis of the accuracy and implications of simple method for predicting the secondary structure of globular proteins）（以下、GOR法と略す）、1987年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」（Journal of Molecular Biology）の第198巻425-443頁掲載のギブラト（Gibrat）らによる論文「ファザー デベロプメンツ オブ プロテイン セコンダリー ストラクチャー プレディクション ユー ジング インホメーション セオリー：ニュー パラメータズ アンド コンシダレーション オブ レジデュアー ペアズ」（Further developments of protein secondary structure prediction using information theory: New parameters and consideration of residue pairs）（以下、GGR法と略す）、及び1988年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」（Journal of Molecular Biology）の第202巻865-884頁掲載のキャン（Qian）らによる論文「プレディクティン グ ザ

セコンダリー ストラクチャー オブ グローブ  
 プロテインズ ユージング ニューラル ネットワーク  
 モデルズ」(Predicting the secondary structure of globular proteins using neural network models) (以下、QS法と略す) などがある。CF法は、タンパク質構造のデータベースから各二次構造におけるアミノ酸の統計的な出現頻度を求め、この頻度表を使用し、経験的な規則に基づく予測を行っている。また、GOR法は、中心残基の二次構造に対して、その残基から数残基離れた残基により独立にもたらされる情報量の和を計算し、その相対値から予測を行い、GGR法は、中心残基の二次構造に対して、その残基及びその残基から数残基離れた残基によりもたらされる情報量の和から予測を行っている。さらに、QS法は、3層のフィードフォワード型のネットワークを使用し、中心残基の前後8残基を含む配列を入力とし、二次構造に対する中心残基及び周辺残基からの寄与をニューラルネットワークを用いて抽出することにより予測を行っている。

【0003】

【発明が解決しようとする課題】 3種類の二次構造の中からアミノ酸配列の各残基に対応する二次構造を選択する予測を3状態予測と呼ぶが、その予測結果の尺度である予測率は、従来の技術のいずれの方法も3状態予測で60%台であり、 $\alpha$ ヘリックスにだけ限ってより予測率の高い予測手法が望まれていた。また、従来の予測結果は、アミノ酸一次配列内の各中心残基に対応する二次構造を予測する残基対応の予測であり、一次配列内のどの領域がどの二次構造に相当するかといった領域対応の予測を行うことも重要であるにも関わらず、このような予測方式に十分な検討がなされていなかった。さらに、アミノ酸配列を文字列としてのみならず、そのアミノ酸の性質(疎水性、分子量など)を考慮した予測を行うことによる予測法も全く確立されていなかった。

【0004】

【課題を解決するための手段】 第1の発明は、タンパク質のアミノ酸配列からタンパク質の構造予測を行うための訓練データを抽出する訓練データ抽出手段と、訓練データから確率的規則の学習を行う学習手段と、学習された確率的規則を用いてテストアミノ酸配列データに対してテストを行うテスト手段とから成ることを特徴とする。

【0005】 第2の発明は、前記訓練データ抽出手段が、構造既知のタンパク質のアミノ酸配列に対して、同じファミリーに属するタンパク質のアライメント(整合)をとり、予測対象とする二次構造領域に対応する部分配列を、二次構造領域の正例として抽出するステップと、構造既知のタンパク質の予測対象とする二次構造に対応する部分配列に対して、構造既知のタンパク質から

なるデータベースの各配列のアライメントをとり、予測対象とする二次構造に対応しない部分配列を、二次構造領域の負例として抽出するステップとから成ることを特徴とする。

【0006】 第3の発明は、前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の種類から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定することを特徴とする。

【0007】 第4の発明は、前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の実数値属性から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定することを特徴とする。

【0008】 第5の発明は、前記学習手段が、前記正例と前記負例とからなる学習データのアミノ酸の実数値属性から、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定するステップと、前記ステップの確率的規則におけるモデルを情報量規準を用いて最適化するステップとから成ることを特徴とする。

【0009】 第6の発明は、前記テストデータ配列に対するテスト手段が、前記学習方法により学習された確率的規則を使用し、テストデータ配列の各領域に対して、その活動度を計算するステップと、計算された活動度の中から最適値を選出するステップとから成ることを特徴とする。

【0010】

【実施例】 次に、本発明について図面を参照して詳細に説明する。

【0011】 図1は、本発明のタンパク質立体構造予測方法の実施例を説明するフローチャートである。本実施例では、対象とする二次構造として $\alpha$ ヘリックスを扱うものとする。

【0012】 ステップ10は、第2の発明に含まれる。このステップでは、 $\alpha$ ヘリックスの領域がわかっているタンパク質のアミノ酸配列に対して、同じファミリーのタンパク質、例えば、種が異なる同じタンパク質のアライメント(整合)をとり、 $\alpha$ ヘリックスに対応する部分配列を、 $\alpha$ ヘリックスの正例として抽出する。

【0013】 例えば、ヘモグロビンというタンパク質の $\beta$ 鎖の場合には、ヒトのヘモグロビンの $\alpha$ ヘリックスの位置は、X線結晶回折の結果から明らかになっており、8個の $\alpha$ ヘリックスの領域を有することが知られている。従って、ヒトのヘモグロビン $\beta$ 鎖に対して、他の種、例えば、チンパンジー、ウマなどの他の種のヘモグロビン $\beta$ 鎖のアライメントをとり、8個の $\alpha$ ヘリックスに対応する領域を $\alpha$ ヘリックスの正例として抽出する。

【0014】 ステップ20は、第2の発明に含まれる。このステップでは、 $\alpha$ ヘリックス位置の知られているタンパク質の $\alpha$ ヘリックスに対応する部分配列に対して、 $\alpha$ ヘリックス位置の知られているアミノ酸配列データベースの各配列のアライメントをとり、 $\alpha$ ヘリックスに対

応しない部分配列を、ステップ10で抽出された $\alpha$ ヘリックスの正例に対する負例として抽出する。

【0015】ヘモグロビン $\beta$ 鎖の例では、8個の $\alpha$ ヘリックスに対応する部分配列に対して、例えば、PDB (Protein Data Bank) などのタンパク質構造データベース内のいくつかの配列に対してアライメントを行い、アライメントの結果得られた各部分配列において、その配列の構造が $\alpha$ ヘリックスではない場合に、それらを負例として抽出する。例えば、負例抽出の際のアライメントでは、一定の割合以上の相同性を保持する部分配列を負例とすることが考えられる。具体的には、アライメントによる相同性が30%以上の部分配列を負例とする方法などがある。

【0016】抽出するデータ数については、例えば、 $\alpha$ ヘリックスの正例となる各領域における正例と負例との割合を各領域についてそれぞれ等しくすることが考えられ、また例えば、その割合として正例、負例を同数とすることが考えられる。

【0017】ステップ30は、第3の発明、第4の発明、第5の発明に共通に含まれ、確率的規則の実数値パラメータを推定するステップである。このステップでは、ステップ10で求めた正例とステップ20で求めた負例からなる学習データから、確率的規則を用いることにより、この確率的規則の実数値パラメータを推定する。このステップでの確率的規則の構造を、以下に示す。

【0018】確率的規則とは、ここでは任意の与えられた配列の領域に対して、 $\alpha$ ヘリックスが対応する確率を与える確率分布のことである。各 $x_i$  ( $i=1, \dots, n$ ) をそれぞれ属性値の空間として、 $x$ をそれらの直積、すなわち、 $x = x_1 \times x_2 \times \dots \times x_n$  と書く。

【0019】例えば、 $x$ は20種類のアミノ酸からなる一つの集合を表す場合や、また $x = x_1 \times x_2$  で、 $x_1$  が疎水性を表す数値の範囲かつ $x_2$  が分子量を表す数値の範囲を表す場合などがある。この例での前者の場合が第3の発明で使用され、それ以外の場合が第4の発明及び第5の発明で使用される。Sをある領域の長さWの配列であり、各Sは $x \times x \times \dots \times x$  の元とみなし、ま\*

\*た、 $X_i$  を配列Sの左から数えてi番目の残基であり、 $P(\alpha | X_i)$  が、 $X_i$  に対応する二次構造が $\alpha$ ヘリックスである確率とする。ここで、配列Sに対応する二次構造が $\alpha$ ヘリックスである確率 $P(\alpha | S)$  は、 $P(\alpha | X_i)$  の積として次のようにかけるものと仮定する。

【0020】

$$P(\alpha | S) = \prod_{i=1}^W P(\alpha | X_i)$$

さらに、各 $P(\alpha | X_i)$  の具体的表現として、例えば、有限分割型確率的規則を使用する。有限分割型確率的規則は次のような構造をもつ条件付き確率分布であり、以下のように構成する。前記配列Sのi番目の残基における属性の実数値のとり得る範囲を重なり合わない部分領域（以下、これをセルと呼ぶ）に分割し、mを全セル数、 $C_k$  をk番目のセルとした時に、 $X_i$  がm個のセルの内の $C_k$  に含まれる場合に、 $P(\alpha | X_i) = P_k(i)$  とする。ここで、

【0021】

【数1】

$$P_k(i) \in [0, 1] \quad (k=1, \dots, m)$$

【0022】であり、これを確率パラメータと呼ぶ。図4は、有限分割型確率規則の構造を示す模式図であるが、この図では、一例として、値が0から1の範囲をとる一つの属性により確率パラメータを推定する場合を示す。

【0023】確率パラメータは、各セルに含まれる正例及び負例のデータ数を用いて推定する。mをセルの数、 $N_k^+(i)$  をi番目の位置でのk番目のセルに含まれる正例数、 $N_k^-(i)$  をi番目の位置でのk番目のセルに含まれる負例数、 $N_k(i)$  をi番目の位置でのk番目のセルに含まれる正例数と負例数の和とし、i番目の位置でのk番目のセルにおける推定値を

【0024】

【数2】

$$\hat{P}_k(i)$$

【0025】とする。例えば、次式のラプラス推定量によって、各セルに対する確率パラメータを計算する。

【0026】

【数3】

$$\hat{P}_k(i) = (N_k^+(i) + 1) / (N_k(i) + 2)$$

$$(k=1, \dots, m)$$

【0027】ただし、推定量はラプラス推定量のみならず、多くの推定量が使用できる。

【0028】ステップ40は、第6の発明に含まれる。このステップでは、ステップ30において学習された確率的規則を使用し、テストデータ配列の各領域に対して、その活性度を計算する。

【0029】ここでは、活性度として尤度を使用する。

【0030】具体的には、確率的規則が構成された長さ

wのある $\alpha$ ヘリックス領域を考える。テストデータの $\alpha$ ヘリックス配列に対して、前記領域の長さwより小さな長さtのw-t+1個のすべての部分領域を設け、このw-t+1個の部分領域それぞれをテスト $\alpha$ ヘリックス配列の左から順にあてはめていき、テスト配列の各領域の尤度を計算する。

【0031】さて、k番目の長さtの部分領域に対して、 $\alpha$ ヘリックス領域の確率パラメータを左から順に並

べたものを  $\xi_k = (\theta_1, \dots, \theta_t)$ ,  $\theta_i = (P_1(i), \dots, P_m(i))$  ( $i=1, \dots, t$ ) と書く。

【0032】ここで、 $m$  はセルの数であり、 $\theta_i$  は既に学習によって値が求められている。

【0033】 $w-t+1$  個の部分領域の位置に対応して、この  $mt$  次元パラメータは、 $w-t+1$  個求められるので、それを  $\xi_1, \dots, \xi_{w-t+1}$  とする。

【0034】前記パラメータを使用して、任意の長さ  $t$  のテストアミノ酸配列  $\Gamma$  に対して、尤度が  $w-t+1$  通り次のように計算できる。

【0035】

$$P(\alpha | \Gamma : \xi_k) \quad (k=1, \dots, w-t+1)$$

ただし、各  $k$  について、

$$P(\alpha | \Gamma : \xi_k) = \prod_{i=1}^t P(\alpha | \Gamma : \theta_i)$$

ここで、 $P(\alpha | \Gamma : \theta_i)$  は  $X_i$  が  $l$  番目のセルに入れば、 $P_l(i)$  ( $l=1, \dots, m$ ) と計算する。また、 $P_l(i)$  ( $l=1, \dots, m$ ) はすでに学習されている。

【0036】例えば、前記有限分割型確率的規則でのアミノ酸の属性値の空間がある一つの属性値のみからなる場合で、またセルの数が3であり、セルに入る推定量はラプラス推定量により求めるとする。このとき、ある  $\alpha$  ヘリックス領域の  $k$  番目の長さ5の部分領域の  $i$  番目の位置での  $l$  番目のセルの正例数を  $N_l^+(i)$ 、 $l$  番目のセルの負例数を  $N_l^-(i)$ 、正例数と負例数の和を  $N_l(i)$  とする。すると、 $i$  番目の位置での  $l$  番目のセルの推定量は、例えば、 $P_l(i) = (N_l^+(i) + 1) / (N_l(i) + 2)$  として得られる。

【0037】ここで、テストアミノ酸配列のウィンドウの大きさ5の領域  $\Gamma$  に対してテストを行うとし、領域  $\Gamma$  のそれぞれの残基は前記部分領域での構成された学習規則での各1, 2, 3, 2, 1番目のセルに入る属性の実数値を有するとする。すると、前記  $k$  番目の部分領域によるテストアミノ酸配列の領域  $\Gamma$  の活性度は尤度  $P(\alpha | \Gamma : k)$  として、次式のように計算される。

$$P(\alpha | \Gamma : k) = \{ (N_1^+(1) + 1) / (N_1(1) + 2) \} \{ (N_2^+(2) + 1) / (N_2(2) + 2) \} \{ (N_3^+(3) + 1) / (N_3(3) + 2) \} \{ (N_2^+(4) + 1) / (N_2(4) + 2) \} \{ (N_1^+(5) + 1) / (N_1(5) + 2) \}$$

$w-t+1$  個の部分領域により、テストアミノ酸配列のとり得るすべての領域に対して、この尤度計算を行う。また、 $\alpha$  ヘリックス領域が複数個存在すれば、その各領域について同様の尤度計算を行う。

【0038】従って、テストアミノ酸配列内でのウィンドウの大きさに対応したすべての領域に対して、尤度が出力として得られることになる。

【0039】以上のウィンドウを使用した各領域に対応

する尤度計算により、一つ一つの残基に対して  $\alpha$  ヘリックスが対応する確率を計算するのではなく、テストアミノ酸配列の各部分領域に  $\alpha$  ヘリックスが対応する確率を尤度として計算することができる。

【0040】ステップ50は、第6の発明に含まれる。このステップでは、前記ステップ40により計算された複数の活性度の中で、 $\Gamma$  に対して最適な一つの活性度を求め、さらにテストアミノ酸配列全体における活性度の変化を出力する。

【0041】ステップ40に引き続きここでは、活性度として尤度を使用する。

【0042】例えば、最適値  $P(\alpha | \Gamma : \xi_k^*)$  を以下で定める。

$$P(\alpha | \Gamma : \xi_k^*) = \max \{ P(\alpha | \Gamma : \xi_1), \dots, P(\alpha | \Gamma : \xi_{w-t+1}) \}$$

$\alpha$  ヘリックス領域が複数個あれば、各領域について、同じ  $\Gamma$  に対して同様の尤度計算を行ない、 $\alpha$  ヘリックス領域全体を通じて最大の尤度を最適値として選ぶ方法も考えられる。

【0043】さらに、テスト配列内の尤度が与えられた各領域において、最大の尤度を領域内の各残基の最適値とする、あるいは、領域内の各残基に対しては、その残基を含む領域の得られた尤度の平均を各残基の最適値とする、などの方法を使用し、テストアミノ酸配列全体に対する尤度の変化を出力する。

【0044】以上の図1における学習及び予測方法は、 $\alpha$  ヘリックス以外の二次構造予測についても適用できる。

【0045】図2は、本発明のタンパク質立体構造予測方法の実施例を説明するフローチャートである。本実施例では、対象とする二次構造として  $\alpha$  ヘリックスを扱うものとする。

【0046】ステップ60は、図1のステップ10と同じ処理を行ない  $\alpha$  ヘリックス領域予測のために必要な正例を抽出する。

【0047】ステップ70は、図1のステップ20と同じ処理を行ない  $\alpha$  ヘリックス領域予測のために必要な負例を抽出する。

【0048】ステップ80は、図1のステップ30と同じ処理を行ない確率的規則を適用し、この確率的規則の実数値パラメータを推定する。

【0049】ステップ90は、第5の発明にのみ含まれる。このステップでは、確率的規則のモデルを情報量規準を用いて最適化する。使用する情報量規準としては、例えば、MDL (minimum description length) 規準などが考えられる。

【0050】前記有限分割型確率的規則にあてはめると、MDL原理は、データ記述長と有限分割型確率的規則による記述長との和が最小である時に最適な確率的規則が構成されているとする。なお、MDL原理について

は、1978年発行の米国の雑誌「オートマティカ」(Automatica)の第14巻465-471頁記載のリサネン(Rissanen)による論文「モデリングバイショーテストデータディスクリプション」(Modeling by shortest data description)に詳しく記載されている。

\*

$$-\sum_i \sum_k \{N_k^+(i) \log p_k(i) + (N_k(i) - N_k^+(i)) \log (1 - p_k(i))\}$$

【0053】また、前記有限分割型確率的規則の記述長は、各確率パラメータ $p_k(i)$ の推定値はおおよそ $\log N_k(i)$ ビットで記述できるから、次のように計算できる。

【0054】

【数5】

※

$$\sum_i \sum_k (\log N_k(i)) / 2.$$

【0055】したがって、MDL原理によれば、次式が最小になるセル数 $m$ の大きさを確率規則を構成する最適なセル数とする。

【0056】

※

【数6】

$$-\sum_i \sum_k \{N_k^+(i) \log p_k(i) + (N_k(i) - N_k^+(i)) \log (1 - p_k(i))\} + \sum_i \sum_k (\log N_k(i)) / 2.$$

【0057】ステップ100は、図1のステップ40と同じ処理を行ない、ステップ90を使用してモデルが最適化された確率的規則を使用し、テストアミノ酸配列データの各領域に対して、その活性度を計算する。

【0058】ステップ110は、図1のステップ50と同じ処理を行ないステップ40により求められた複数の活性度から、配列全体に対する活性度の変化を出力する。

【0059】以上の図2における学習及び予測方法は、 $\alpha$ ヘリックス以外の二次構造予測についても適用できる。

【0060】

【発明の効果】二次構造既知のタンパク質のアミノ酸配列情報から、二次構造未知のタンパク質の二次構造を従来技術に対して高い精度で予測することができる。特に、アルブミンの $\alpha$ ヘリックス領域を70%以上の高い精度で予測可能である。また、MDL原理などの情報量規準によりモデルを最適化することにより、確率的規則の構造を理論的に最適化することが可能になる。

【図面の簡単な説明】

【図1】本発明のタンパク質立体構造予測方法の一実施例を示すフローチャート

【図2】本発明のタンパク質立体構造予測方法の一実施例を示すフローチャート

【図3】本発明の二次構造( $\alpha$ ヘリックス)領域予測方法の模式図

【図4】本発明で使用する確率規則の一例である有限分割型確率規則の具体例を示す模式図

30 【符号の説明】

10 正例抽出

20 負例抽出

30 確率的規則による実数値パラメータ推定

40 テスト配列各領域に対する活性度計算

50 テスト配列に対する予測値算出

60 正例抽出

70 負例抽出

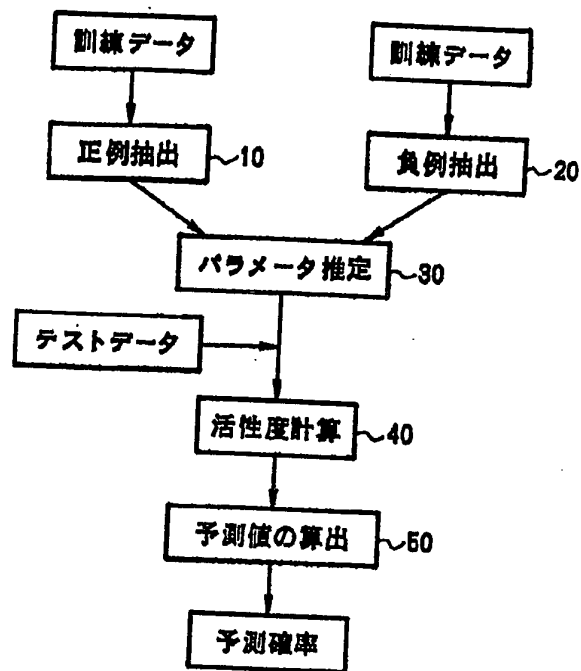
80 確率的規則による実数値パラメータ推定

90 情報量規準による最適化

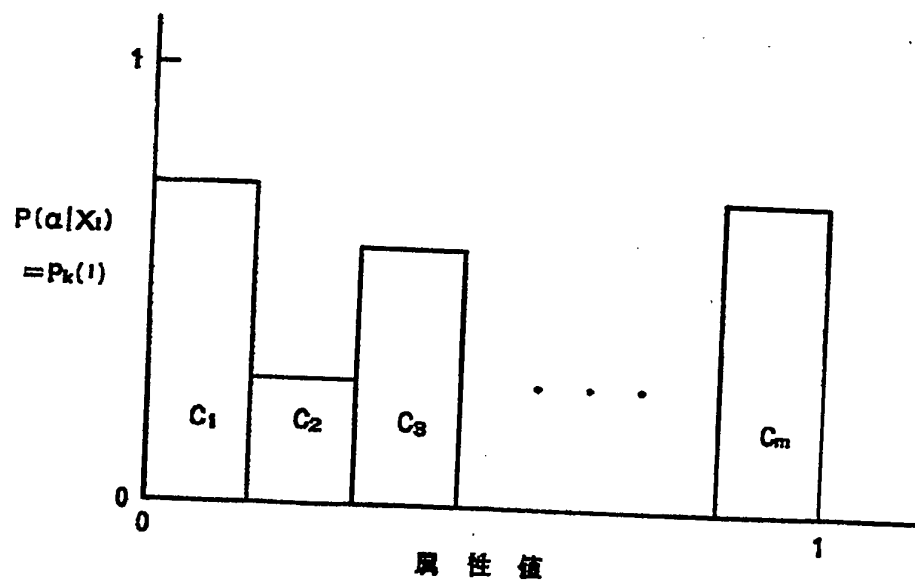
40 100 テスト配列各領域に対する活性度計算

110 テスト配列に対する予測値算出

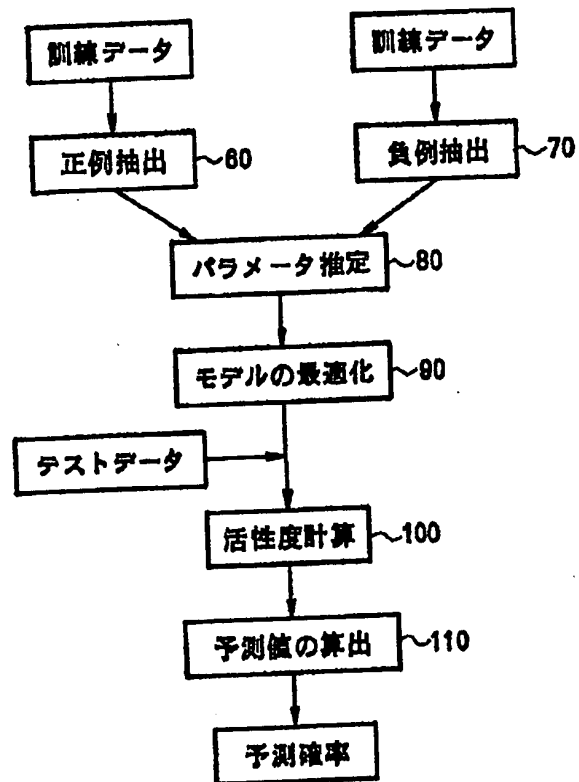
【図1】



【図4】



【図2】





(9)

【図3】

